

1 **ABSTRACT**

2

3 Described herein is a technology for recognizing the content of text

4 documents. The technology determines one or more hash values for the content of

5 a text document. Alternatively, the technology may generate a “sifted text” version

6 of a document. In one implementation described herein, document recognition is

7 used to determine whether the content of one document is copied (i.e., plagiarized)

8 from another document. This is done by comparing hash values of documents (or

9 alternatively their sifted text). In another implementation described herein,

10 document recognition is used to categorize the content of a document so that it

11 may be grouped with other documents in the same category. This abstract itself is

12 not intended to limit the scope of this patent. The scope of the present invention is

13 pointed out in the appending claims.

14

15

16

17

18

19

20

21

22

23

24

25